

Open Research Online

The Open University's repository of research publications and other research outputs

Reducing the effort for systematic reviews in software engineering

Journal Item

How to cite:

Osborne, Francesco; Muccini, Henry; Lago, Patricia and Motta, Enrico (2019). Reducing the effort for systematic reviews in software engineering. *Data Science*, 2(1-2) pp. 311–340.

For guidance on citations see [FAQs](#).

© 2019 IOS Press and the authors



<https://creativecommons.org/licenses/by/4.0/>

Version: Version of Record

Link(s) to article on publisher's website:
<http://dx.doi.org/doi:10.3233/DS-190019>

Copyright and Moral Rights for the articles on this site are retained by the individual authors and/or other copyright owners. For more information on Open Research Online's data [policy](#) on reuse of materials please consult the policies page.

oro.open.ac.uk

Reducing the effort for systematic reviews in software engineering

Francesco Osborne^{a,*}, Henry Muccini^b, Patricia Lago^c and Enrico Motta^d

^a Knowledge Media Institute, The Open University, UK

E-mail: francesco.osborne@open.ac.uk; ORCID: <https://orcid.org/0000-0001-6557-3131>

^b DISIM Department, University of L'Aquila, Italy

E-mail: henry.muccini@univaq.it; ORCID: <https://orcid.org/0000-0001-6365-6515>

^c Department of Computer Science, Vrije Universiteit Amsterdam, The Netherlands

E-mail: p.lago@vu.nl; ORCID: <https://orcid.org/0000-0002-2234-0845>

^d Knowledge Media Institute, The Open University, UK

E-mail: enrico.motta@open.ac.uk; ORCID: <https://orcid.org/0000-0003-0015-1952>

Editor: James McCusker (<https://orcid.org/0000-0003-1085-6059>)

Solicited reviews: Anita Dewaard (<https://orcid.org/0000-0002-9034-4119>); one anonymous reviewer; Philipp Cimiano (<https://orcid.org/0000-0002-4771-441X>)

Received 7 January 2019

Accepted 7 June 2019

Abstract.

CONTEXT: Systematic Reviews (SRs) are means for collecting and synthesizing evidence from the identification and analysis of relevant studies from multiple sources. To this aim, they use a well-defined methodology meant to mitigate the risks of biases and ensure repeatability for later updates. SRs, however, involve significant effort.

GOAL: The goal of this paper is to introduce a novel methodology that reduces the amount of manual tedious tasks involved in SRs while taking advantage of the value provided by human expertise.

METHOD: Starting from current methodologies for SRs, we replaced the steps of keywording and data extraction with an automatic methodology for generating a domain ontology and classifying the primary studies. This methodology has been applied in the Software Engineering sub-area of Software Architecture and evaluated by human annotators.

RESULTS: The result is a novel Expert-Driven Automatic Methodology, EDAM, for assisting researchers in performing SRs. EDAM combines ontology-learning techniques and semantic technologies with the human-in-the-loop. The first (thanks to automation) fosters scalability, objectivity, reproducibility and granularity of the studies; the second allows tailoring to the specific focus of the study at hand and knowledge reuse from domain experts. We evaluated EDAM on the field of Software Architecture against six senior researchers. As a result, we found that the performance of the senior researchers in classifying papers was not statistically significantly different from EDAM.

CONCLUSIONS: Thanks to automation of the less-creative steps in SRs, our methodology allows researchers to skip the tedious tasks of keywording and manually classifying primary studies, thus freeing effort for the analysis and the discussion.

Keywords: Systematic reviews, software engineering, ontology learning, semantic web, software architecture, digital libraries

*Corresponding author. E-mail: francesco.osborne@open.ac.uk.

1. Introduction

Understanding the state-of-the-art in research provides the foundation for building novelty. In particular, in Software Engineering topic areas, the acquisition of knowledge for this understanding follows a clear path: started with informal reviews and surveys, it is moving towards systematic searches of the literature. Kitchenham [21] clearly explains the reasons, the importance, and the advantages and disadvantages of using systematic reviews instead of informal ones. Various studies (e.g., [12,68]) reveal the growing interest in systematic literature reviews and systematic mapping studies [65]. A number of articles and books have been written on how to perform such systematic studies [23,61,64].

A Systematic Review (SR) is “*a means of evaluating and interpreting all available research relevant to a particular research or topic area or phenomenon of interest*” [21]. Given a set of research questions, and by following a systematically defined and reproducible process, a SR helps selecting primary studies that contribute to provide an answer to them. Used in combination with keywording [41], a SR supports the systematic elicitation of an ontological classification framework [42]. In this paper we focus specifically on the field of Software Engineering, but systematic literature reviews and mapping studies are used in several research fields, such as Biomedics [10], Robotics [5], Artificial Intelligence [45], Human-Computer Interaction [30], Psychology [46], and many others.

A SR can help researchers and practitioners in creating a complete, comprehensive and valid picture of the state-of-the-art about a given theme when the search-space is bounded (e.g., when the search query returns few thousands of articles to scrutinize). However, it falls short when used to investigate the state-of-the-art on an entire research area (e.g., Software Architecture) where the returned entries are hundreds of thousands – hence clearly unmanageable. As reported by Vale et al. [60] while investigating the state-of-the-art of the Component-based Software Engineering area through an SR, a “*...manual search [restricted only to the most relevant journals and conferences related to the CBSE area] was considered as the primary source, given the infeasibility of analyzing all studies collected from automatic search*”. Still, they had to select, read, and thoroughly analyze 1,231 primary studies.

In contrast to manually run SRs, several state of the art automated methods allow classifying a document in a certain category or topic [2,8,32,57]. Unfortunately, most current techniques suffer from limitations that make them unsuitable for systematic reviews. The approaches which exploit keywords as proxy for research areas are unsatisfactory, as they fail to distinguish research topics from other terms that can be used to annotate papers (e.g., “user case”, “scalability”) and to take advantage of the relationships that hold between research areas (e.g., the fact that “Software Architecture” is a sub-area of “Software Engineering”). Probabilistic topic models (e.g., Latent Dirichlet Allocation [8]) are also unsuitable for this task since they produce cluster of terms that are not easy to map to research areas [38]. Crucially, it is often unfeasible to integrate these topic detection techniques with the needs and the knowledge of human experts. Another alternative is to apply entity linking techniques [32] to map papers to relevant entities in knowledge base. Unfortunately, we currently lack good granular and machine readable representation of research areas in many domains which could be used to this end.

Current techniques have complementary limitations when investigating the state-of-the-art of an entire research area: on the one hand side, SRs are “*human-intensive*”, as they require domain experts to invest a large amount of time to carry out manual tasks; on the other side, automated techniques keep the humans “*out of the loop*”, while human expertise is critical for the more conceptual analysis tasks.

This paper proposes an *expert-driven automatic methodology* (EDAM) for assisting systematic reviews that, while recognizing the essential value of human expertise, limits the amount of tedious tasks the expert has to carry out. Our methodology contributes with (1) automatically extracting an ontology

of relevant topics, related to a given research area; (2) using experts to refine this knowledge base; (3) exploiting this knowledge base for classifying relevant papers that may be then further validated/analyzed by experts, and for computing research analytics. We demonstrate EDAM in the field of Software Architecture, but it can be easily applied to other research areas as well.

Naturally, the ability of domain experts to analyse the research dynamics emerging from primary studies and to distill the most important lessons and trends is still crucial. Therefore, our aim is not to fully automatize the process, but to assist domain experts by automatically generating data-driven analytics in order to free time and resources for the analysis phase.

In summary, our contributions are:

- a novel methodology for supporting ontology-driven systematic reviews, which involves both automatic techniques and human experts;
- an implementation of this methodology which exploits the Klink-2 algorithm for generating the domain ontology in the field of Software Architecture;
- an illustrative analysis of the Software Architecture trends;
- an evaluation involving six human annotators, which shows that the classification of primary studies yielded by the proposed methodology is comparable to the one produced by domain experts ($p = 0.77$).
- an automatically generated ontology of Software Engineering, which could support further systematic reviews in the field.¹

The rest of the paper is structured as follows. Section 2 introduces related works on systematic studies. Section 3 provides an overview of some preliminary evidence of the benefits brought by using EDAM to assist a mapping study. Section 4 then presents the EDAM methodology and its application to the research area of Software Architecture. This experiment is discussed and evaluated in Section 5, which also present a comparison of several approaches for classifying research papers. Finally, in Section 6 we discuss the main implications of our study and outline future directions of research.

2. Related work

There are many guidelines for, and reports on, carrying out systematic studies in Software Engineering. Among them, we could identify a few aimed at supporting or improving the underlying process. In our perspective, they all enable researchers to focus more on the most creative steps of a systematic study by removing what is referred to as *manual work*.

With a motivation similar to ours, i.e. to improve the search step in systematic studies in Software Engineering research, Octaviano et al. [35] propose a strategy that automates part of the primary study selection activity. Mourão et al. [34] present a preliminary assessment of a hybrid search strategy for systematic literature reviews that combines database search and snowballing to reduce the effort due to searches in multiple digital libraries. Kuhrmann et al. [26] provide recommendations specifically for the general study design, data collection, and study selection procedures. Zhang et al. [69], in turn, systematically select and analyze a large number of SRs. Their results have been then used to define a quasi-gold standard for future studies. In their validation, they were able to improve the rigor of the search process and provide guidelines complementing the ones already in use.

¹<http://rexplore.kmi.open.ac.uk/data/edam/SE-ontology.owl>

Ros et al. [48] propose a machine learning approach that classifies papers for SRs by leveraging human experts, who iteratively validate set of publications produced by a classifier. Conversely, EDAM does not require experts to manually examine research papers, but only to review a taxonomy of research areas.

The need for guidelines in conducting empirical research has been addressed in other types of empirical studies, too. De Mello and Travassos [14] focus on opinion surveys and provide guidelines (in the form of a reference framework) aimed at improving the representativeness of samples. Also on opinion surveys, Moller et al. [33] provide recommendations based on an annotated bibliography instead.

Another interesting work by Felizardo et al. [16] investigates how the use of forward snowballing can considerably reduce the effort in updating SRs in Software Engineering. Based on this result, complementing our method with automated forward snowballing suggests a very promising direction for future works as it could further reduce the effort for identifying relevant primary studies.

Marshall et al. [31] carried out an interview survey with experts in other domains (i.e. healthcare and social sciences) with the aim to identify tools that are generally used, or desirable, to ease which steps in systematic studies, and transfer the best practices to the Software Engineering domain. Among the results, data extraction and automated analysis emerge as top requirements for reducing the workload. In a similar vein, Hassler et al. [20] followed by Al-Zubidy et al. [1] consulted Software Engineering researchers conducting SRs to identify and prioritize the necessary SR tool features. The results identified *search & study selection* as the most desirable feature. Our work addresses the needs identified by both Marshall et al. [31] and Hassler et al. [20].

The idea of using ontologies for supporting SRs was discussed by few papers, but did not receive much attention. de Almeida Biolchini et al. [13] introduced the Scientific Research Ontology, a resource to organize the knowledge generated from SR. This ontology offers a conceptual framework with the aim of fostering the consistency between different studies, but does not directly assist the tasks involved in SR, such as the extraction of primary studies. Sun et al. [59] discussed the use of ontologies for supporting key activities in SRs and presented an experiment in which they automatically classified primary studies by means of COSONT, an ontology of methods for cost estimation. Unfortunately, their approach still required to manually check hundreds of papers and the COSONT ontology was quite simplistic, being an handcrafted list of methods with no hierarchical structure. This is a common issue with manually generated ontology of research concepts, which are usually costly to produce, coarse-grained, and slow to evolve [37]. Conversely, EDAM takes advantage of recent ontology learning techniques to automatically generate complex multi-level ontologies (e.g., the SE ontology presented in this paper includes 956 topics and 5,461 relationships), exploits the resulting taxonomic structure to classify the primary studies, and does not require experts to manually review a large number of papers.

3. An overview of the benefits of automatic SRs

Before entering the details of the EDAM methodology, this section provides an overview of the benefits such an automatic SR methodology can bring with respect to more traditional, manual SRs carried out according to predefined protocols.

We all agree that manual SRs based on well-defined systematic protocols help reducing (but not fully removing) subjective biases in the selection of the studies. They however are by and large unfeasible in reviewing a too large dataset (i.e. when the number of scientific publications is too large to be manually processed by the researcher).

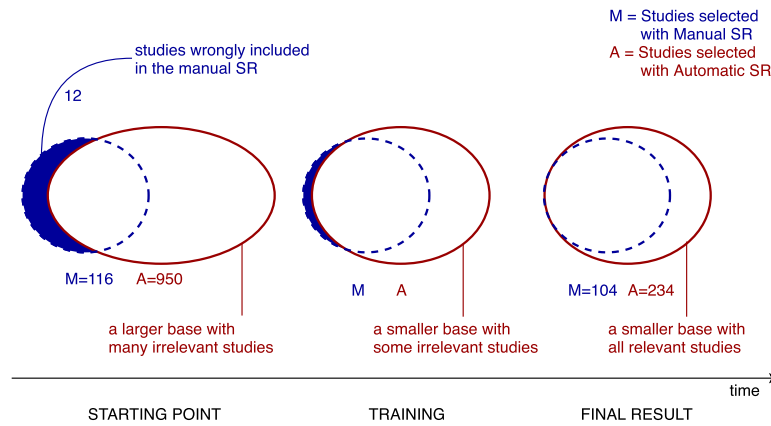


Fig. 1. Some evidence on the benefits of automated SRs.

In a similar vein, automatic SRs help reducing subjective biases (in this case by *implementing* the selection of the studies according to the predefined systematic protocol). Differently, they pose no limitation in terms of the size of the dataset of publications.

In our earlier work [67] we challenged these limitations and benefits by applying the automatic study selection to a manual SR carried out beforehand by other researchers [17]. In this way, we could compare and contrast the results of the manual SR with the results of our automatic SR. In this earlier work, we have studied the field of software sustainability within the Software Engineering domain. While at the time the EDAM methodology was not yet fully developed, we did use the same ontology-learning algorithms and a preliminary version of the ontology for the Software Engineering domain.

The observations gathered during this experiment are illustrated in Fig. 1, where we represented the primary studies selected manually (see the left-hand circles) and those selected automatically (see the right-hand ovals). The experiment underwent three phases:

Starting point: The already-completed manual SR had selected 116 primary studies. Before training the algorithm and tuning the domain ontology, from the Scopus dump of scientific publications we automatically selected 950 studies. While our automatic methodology is able to handle seamlessly any size of the base of publications, the selected studies did initially include a very large number of false positives. However, it did also uncover that 12 studies selected in the manual SR were wrongly included. **Observation #1:** *in spite of systematic selection criteria and the involvement of multiple researchers, human errors in the manual study selection is still possible.*

Training: By treating the 104 primary studies (from the manual SR) as pilot studies, we trained our domain ontology and learning algorithm to automatically select the primary studies. **Observation #2:** *Automatic SR is able to automatize the selection criteria of systematic reviews while handling any size of the initial dataset of scientific publications.* As discussed in Section 5.1, the domain ontology is able to classify the primary studies as correctly as the human experts do, without needing further training. As such, the domain ontology can be reused for any study in the domain of Software Engineering.

Final result: The final result of the automatic selection converged to 234 studies which included the 104 pilot studies and *correctly* identified additional 130 studies that were missing in the original manual SR. **Observation #3:** *By handling a much larger base of publications, automatic SRs are able to uncover primary studies that are missed by manual SRs where such scale is unfeasible.*

4. An expert-driven automatic methodology

We propose a novel expert-driven automatic methodology (EDAM) for assisting systematic reviews like systematic literature reviews and mapping studies. EDAM allows to automatize the steps that are the most time and effort consuming while requiring the least creativity, such as *selection of relevant papers*, *keywording*, and *creation of a classification schema* [42], by exploiting ontology learning techniques and semantic technologies to foster scalability, objectivity, reproducibility, and granularity of the study (further discussed in Section 5.4). It also supports the generation of research trends, which are typical of data synthesis in mapping studies. In this paper, we illustrate how EDAM can support mapping studies, even though it can be evidently exploited in systematic literature reviews, too.

Figure 2 shows the steps of a mapping study using EDAM in contrast with the steps of a classic (manual) methodology – shown in Fig. 3. The main difference is that in the classic methodology the researchers first select and analyze each primary study (steps 2–3) and then produce a taxonomy to classify them (step 4). When assisted by EDAM, instead, the researchers first use ontology learning methods over large scholarly datasets to generate an ontology of the field (steps 2–3), then refine the ontology with the help of domain experts (step 4), and finally exploit this knowledge base to automatically select and classify the primary studies (steps 5–6).

An alternative solution for steps 2–4 (Generation of domain ontology) is the reuse of an ontology crafted by a previous study with the same scope. Indeed, in the study discussed in Section 4.2 we have generated an ontology of Software Engineering (SE) research topics, with the hope that it will be re-used by the research community.

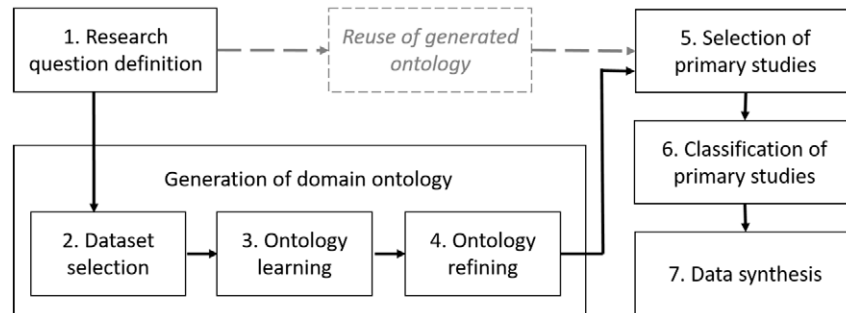


Fig. 2. Steps of a systematic mappings adopting the EDAM methodology. The gray-shaded elements refer to the alternative step of reusing the previously generated ontology.

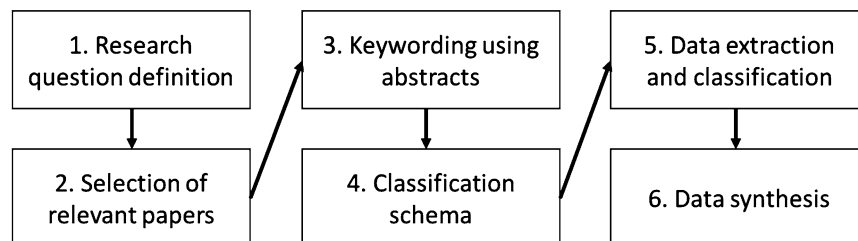


Fig. 3. Classic steps of systematic mappings (inspired by [42]).

In Section 4.1, we describe EDAM and discuss its advantages over a classic methodology. In Section 4.2, we exemplify the application of EDAM specifically aimed at identifying publication trends of the Software Architecture research area in the specific SE domain.

4.1. EDAM description

A SR assisted by EDAM is organized along the following steps (ref. Fig. 2).

1. Research question definition. The researchers performing the study state the research questions (RQs). These will affect the aim of the study and thus its steps. It should be noted that EDAM is applicable only to research questions that could be answered by classifying publications, authors, venues, and other entities according to the ontology for producing relevant analytics. Other research questions should be addressed according to the standard methodology [42].

2. Dataset selection. The researchers select a dataset on which to apply the chosen ontology learning technique (further elaborated in step 3) for generating the domain ontology that will be used to select and classify the primary studies. The most important characteristic of this dataset is that it must be unbiased with respect to the focus of the study. For example, if the study wants to uncover the trends in research areas (e.g., Software Architecture), the dataset should not be biased with respect to any area in the domain (e.g., Software Engineering in our case). A good strategy to select unbiased datasets is considering either a full scholarly dataset of a very high-level field (e.g., all the Computer Science papers in Microsoft Academic Search² or in Scopus³) or a dataset including all the papers published in the main conferences and journals of the domain under analysis. In recent years, universities, organizations, and publishing companies have released an increasing number of open datasets that could assist in this task, such as CrossRef,⁴ SciGraph,⁵ OpenCitations,⁶ DBLP,⁷ Semantic Scholar,⁸ and others.

3. Ontology learning. The dataset is processed by an ontology learning technique that automatically infers an ontology of the relevant concepts.

We strongly advocate the use of an ontology learning technique that generates a full domain ontology and represents it with Semantic Web standards, such as the Web Ontology Language (OWL)⁹. The main advantage of adopting an ontology in this context is that it allows for a more comprehensive representation of the domain since it includes, in addition to hierarchical relationships, also other kinds of relationships (e.g., *sameAs*, *partOf*), which may be critical for classifying the primary studies. For example, an ontology allows to explicitly associate to each category a list of alternative labels or related terms that will be used in the classification phase. In addition, ontology learning techniques can infer very structured multi-level ontologies [37], and thus describe the domain at different levels of granularity.

The task of ontology and taxonomy learning was comprehensively explored over the last 20 years. Therefore, the researcher can choose among a variety of different approaches for this step, including:

- statistical methods for deriving taxonomies from keywords [29,56];
- natural language processing approaches, e.g., FRED [18], LODifier [4], Text2Onto [11];

²<http://academic.research.microsoft.com>

³<https://www.scopus.com/>

⁴<https://www.crossref.org/>

⁵<https://scigraph.springernature.com/explorer/downloads/>

⁶<http://opencitations.net>

⁷<http://dblp.uni-trier.de>

⁸<https://www.semanticscholar.org/>

⁹<https://www.w3.org/OWL/>

- approaches based on deep learning, e.g., recurrent neural networks [43];
- hybrid ontology learning frameworks [63];
- specific approaches for generating research topic ontologies, e.g., Klink-2 [37].

However, as discussed in the following step, researchers may also choose to skip this step and re-use a compatible ontology from a previous study.

It is useful to clarify why we suggest the adoption of an ontology learning approach, rather than the adoption of one of the currently available research taxonomies, such as the ACM computing classification system,¹⁰ the Springer Nature classification,¹¹ Scopus subject areas,¹² and the Microsoft Academic Search classification. Unfortunately, these taxonomies suffer from some common issues, which make them unfeasible to support most kinds of SRs. First, they are very coarse-grained and represent wide categories of approaches, rather than the fine-grained topics addressed by researchers [36]. Secondly, they are usually obsolete since they are seldom updated. For example, the 2012 version of the ACM classification was finalized fourteen years after the previous version. This is a critical point, since some interesting trends could be associated with recently emerged topics. In third instance, most ontology learning algorithms are not limited to learning research areas, but can be tailored to yield the outputs which are more apt to support a specific analysis.

4. Ontology refining. The ontology resulting from the previous step is corrected and refined by domain experts. During this phase, the experts are allowed to (1) delete an existent category, (2) add a new category, (3) delete an existent relationship, (4) add a new relationship. We suggest using at least three domain experts for addressing possible disagreements.

This step is critical for two reasons. First, it may correct some errors in the automatically-generated taxonomy. Secondly, it verifies that the data-driven representation aligns with the domain experts mental model and thus the outcomes will be understandable and reusable by their research community.

Refining a very large ontology is not a trivial task, therefore if the domain comprehends a large number of topics we suggest splitting it in manageable sub branches to be addressed by different experts. Our experience suggests that a taxonomy of about 50 research areas can be reviewed in about 15–30 minutes by an expert of the field. For example, in [37] three experts reviewed a Semantic Web ontology of 58 topics in about 20 minutes. In the test study for this paper, three experts took about 20 minutes to examine and produce feedback on a taxonomy of 46 topics (and 71 terms considering synonymous such as “product line”, “product-lines”, “product-line”, which were clustered automatically by the ontology learning algorithm). In both cases, we represented the ontology as tree diagram in a excel sheet¹³ and included also a list of the most popular terms in the dataset, for supporting experts in remembering all the relevant research topics.

An alternative solution is to provide experts with ontology editors that could be used to directly modify the ontology, such as Protege,¹⁴ NeOn Toolkit,¹⁵ TopBraid Composer,¹⁶ Semantic Turkey,¹⁷ or Fluent Editor.¹⁸ However, these tools are not always easy to learn and the adoption of a simple spreadsheet may

¹⁰<http://www.acm.org/publications/class-2012>

¹¹<http://www.nature.com/subjects>

¹²<https://www.elsevier.com/solutions/scopus/content>

¹³See an example at <http://tinyurl.com/yal6h3wu>.

¹⁴<http://protege.stanford.edu>

¹⁵<http://neon-toolkit.org/>

¹⁶http://www.topquadrant.com/products/TB_Composer.html

¹⁷<http://semanticturkey.uniroma2.it/>

¹⁸<http://www.cognitum.eu/Semantics/FluentEditor/>

be advisable in most cases. Indeed, the annotators who participated in the mapping study of Software Architecture described in the next section, reported that they were able to easily correct and suggest changes in the ontology using this simple solution. In particular, this task was natural to them since the same kind of spreadsheet is typically used in the analysis phase of systematic reviews (e.g., for the keywording step). We refer the reader to Sabou et al. [49] for a comprehensive analysis of the verification of domain knowledge by human experts in the field of Software Engineering.

As highlighted by Fig. 2, the aim of steps 2–4 is to generate an ontology apt to select and classify relevant papers and ultimately answer the RQs. It follows that these steps could be replaced by the adoption of an ontology previously generated and validated by a previous study with a consistent scope. For example, the ontology about Software Engineering generated for this paper's example study (see Section 4.2) can be re-used to perform many kinds of mapping studies involving other research areas in SE. Naturally, the ontology may have to be further updated to include the most recent concepts and terms. This solution allows users with no access to vast scholarly databases or no expertise in ontology learning techniques to easily implement an EDAM study.

5. Selection of primary studies. The authors select a dataset of papers and define the inclusion criteria of the primary studies according to the domain ontology and other metadata of the papers (e.g., year, venue, language). The inclusion criteria are typically expressed as a search string, which uses simple logic constructs, such as AND, OR, and NOT [3]. The search string is then used to produce the query that will be run over the dataset for selecting the primary studies. Some examples of queries include (1) “all the papers in the dataset published in a list of relevant conferences” or “all the papers in the dataset that contain a list of relevant terms from the ontology”.

In most cases this dataset will be the same or a subset of the one used for learning the domain ontology. However, the authors may want to zoom on a particular set of articles, such as the ones published in the main venues of a field, in a geographical area, or by a certain demography. It is also possible to select a different dataset altogether, since the ontology would use generic topic labels and thus be agnostic with respect to the dataset. A possible reason to do so is the availability of the full text of the studies. Many ontology learning algorithms can be run on massive metadata dataset (e.g., Scopus, Microsoft Academic Search), but some research questions may require the full text. In this case, the author may want to perform the ontology learning step on the metadata dataset, which is usually larger in size and scope, and then either select a subset composed by publications which are available online or adopt for this phase a second dataset that includes the full text of the articles, such as Core [24]. The growth of the Open Access movement [62], which aims at providing free access to academic work, may alleviate this limitation in the following years.

6. Classification of primary studies. The authors define a function for mapping categories to papers based on the refined ontology. This step is important to foster reproducibility since the inclusion criteria (defined in the step 5), the mapping function, and the domain ontology should contain all the information needed for replicating the classification process. The function can also be associated to an algorithmic method (e.g., a machine learning classifier), provided the method is made available and is reproducible.

The simplest way for mapping categories to papers is to associate to each category each paper that contains the label of the category or of any of its sub-categories. This simple technique for semantically characterizing documents has the advantage of being unsupervised and was applied with good results in a variety of fields, such as topic forecasting [52], automatic classification of proceeding books [39], sentiment analysis [50], recommender systems [15], and many others. Alternative unsupervised methods, which we evaluate in Section 5.2, include approaches based on TF-IDF [44], LDA [8], and word embeddings [53].

In addition, the authors can choose to create a more complex mapping function which exploits other semantic relationships in the ontology (e.g., *relatedTerm*, *partOf*).

7. Data synthesis. According to the RQs, this step may be automatic, semi-automatic or manual. Some straightforward analytics (e.g., the number of publications or citations over time) can be computed completely automatically by counting the previously classified papers or summing their number of citations. Other more complex analyses may require the use of machine learning techniques or the (manual) intervention of human experts. Starting from the groundwork formed by our research, a full analysis of the possible kinds of data synthesis and the way to automatize them will constitute interesting future works beneficial for the whole research community.

Overall, motivated by the need to reduce the amount of manual tedious tasks involved in SRs, **EDAM offers four main advantages over a classic methodology**. **First**, human experts are not required to manually analyze and classify primary studies, but they simply have to refine the ontology, choose the inclusion criteria, and define a mapping function for associating papers to categories in the ontology. This allows researchers to carry out large scale studies that involve thousands of research papers with relative ease. **Secondly**, since the domain ontology is created with a data-driven method, it should reflect the real trends of the primary studies, rather than arbitrary human decisions about which keywords to annotate and aggregate, even if the refinement step may still introduce a degree of arbitrariness. **Third**, the use of a formal machine-readable ontology language for representing the domain taxonomy should foster the reproducibility of the study and allow authors with no expertise in data science to perform studies using previously generated ontologies. **Fourth**, this methodology allows researchers to produce and exploit complex multi-level ontologies, rather than the simple two-level classifications used by many studies [60].

Naturally, EDAM is suitable for research questions that can be automatized by the ontology-driven classification process previously described, or that aim at giving an overview of the state-of-the-art or state-of-practice on a topic [66] by analysing all of the relevant research contributions in a specific research area. We will discuss further this and other limitations in Section 5.3.

4.2. EDAM application

With the aim of presenting a reproducible pipeline and showing how EDAM can be applied, we present here an example as part of a possible systematic mapping study assisted by EDAM in the Software Architecture research area. We chose to study the research trends in this area, since trend analysis is typical of mapping studies [66] and it is one of the tasks that can be automatized by EDAM.

In the following, we describe how we instantiated the study example assisted by EDAM and discuss the specific technologies used to implement it. The data necessary for reproducing this study and using this same pipeline on other fields are available at <https://doi.org/10.5281/zenodo.2653924>.

1. Research question definition. We wanted to focus on a task that is often addressed by mapping studies and could be completely automatized. Therefore our RQ is: “What are the trends of the main research topics of Software Architecture?”.

2. Dataset selection. We selected all papers in a dump of the Scopus dataset about Computer Science in the period 2005–2013. The Scopus dataset we were given access by Elsevier BV includes papers in 1900–2013 interval, but the number of relevant articles before 2005 was too low to allow a proper trend analysis. Each paper in this dataset is described by title, abstract, keywords, venue, and author list.

3. Ontology learning. We applied the Klink-2 algorithm [37] on the Scopus dump for learning an ontology representing the main ‘Software Architecture’ research area in SE.

Klink-2 is an algorithm that generates an ontology of research topics by processing scholarly metadata (titles, abstracts, keywords, authors, venues) and external sources (e.g., DBpedia, calls for papers, web pages). In particular, Klink-2 periodically produces the Computer Science Ontology (CSO)¹⁹ [54] that is currently used by Springer Nature for classifying proceedings in the field of Computer Science [39], such as the well-known Lecture Notes in Computer Science series.²⁰ The ontologies produced by Klink-2 use the Klink data model,²¹ which is an extension of the BIBO ontology²² that in turn builds upon SKOS.²³ This model includes three semantic relations: *relatedEquivalent*, which indicates that two topics can be treated as equivalent for the purpose of exploring research data; *skos:broaderGeneric*, which indicates that a topic is a subarea of another one; and *contributesTo*, which indicates that the research outputs of one topic significantly contribute to the research into another. In the following, we make use of the first two relationships for classifying studies according to their research topics.

In Algorithm 1, we report the pseudocode of Klink-2. The algorithm takes as input a set of keywords and investigates their relationships with the set of their most co-occurring keywords. Klink-2 infers a sub-topic relationship between keyword x and y by means of two metrics: (i) $H_R(x, y)$, which uses a semantic variation of the subsumption method; (ii) $T_R(x, y)$, which uses temporal information to do the

Algorithm 1: The Klink-2 algorithm

Input : List of keywords *keywords*, Metadata *metadata*

Output: Ontology *ontology*

```

1 relationships={ };
2 while some keywords yet to process do
3   foreach k1 in keywords do
4     candidates = GetCandidates(k1, metadata);
5     foreach k2 in candidates do
6       | relationship = InferRelationship(k1, k2, metadata, relationships);
7     end foreach
8   end foreach
9   relationships = RemoveLoops(relationships);
10  new.keywords = MergeAndSplitKeywords(keywords, metadata, relationships);
11  keywords = AddNewKeywords(new.keywords);
12 end while
13 keywords = FilterTopics(keywords, metadata, relationships);
14 ontology = GenerateSemanticRelationships(relationships);
15 return(ontology);

```

¹⁹<http://cso.kmi.open.ac.uk/>

²⁰<http://www.springer.com/gp/computer-science/lncs>

²¹<http://technologies.kmi.open.ac.uk/rexplorer/ontologies/BiboExtension.owl>

²²<http://purl.org/ontology/bibo/>

²³<https://www.w3.org/2004/02/skos/>

same. $H_R(x, y)$ is computed according to the following formula:

$$H_R(x, y) = \left(\frac{I_R(x, y)}{I_R(x, x)} - \frac{I_R(y, x)}{I_R(y, y)} \right) c_R(x, y) n(x, y)$$

where $I_R(x, y)$ is the number of elements associated with both x and y according to relation R (e.g., number of co-occurrences in research papers), $\frac{I_R(x, y)}{I_R(x, x)}$ is the conditional probability that an element associated with keyword x will be associated also with keyword y , $n_R(x, y)$ is the Levenshtein distance between the two keywords normalized by the length of the longest one, and $c_R(x, y)$ is the cosine similarity between the two vectors in which each index represents a keyword k and the value is the number of time x or y co-occurred with k in a certain context.

$T_R(x, y)$ is a temporal version of $H_R(x, y)$, which weighs more the information associated with the first years of x . It is useful to detect the cases in which the relationship between two terms fade because their association has become implicit (e.g., Artificial Intelligence and Machine Learning). $T_R(x, y)$ is calculated using a variation of formula (1) in which $I_R(x, y)$ is computed by weighting the intensity of the relationships in each year according to the distance from the debut of x . The weight is computed as $w(\text{year}, x) = (\text{year} - \text{debut}(x) + 1)^{-\gamma}$, with $\gamma > 0$ ($\gamma = 2$ in the implementation used for this paper). A hierarchical relationship is inferred whenever $H_R(x, y)$ or $T_R(x, y)$ are higher then a certain threshold (0.25 in the implementation used for this study).

After inferring the hierarchical relationships, Klink-2 removes loops in the topic network (instruction #9), merges similar keywords and splits ambiguous keywords associated to multiple meanings (e.g., ‘Java’). The keywords produced in this step are added to the initial set of keywords to be further analysed in the next iteration and the while-loop is re-executed until there are no more keywords to be processed. Finally, Klink-2 filters the keywords considered ‘too generic’ or ‘not academic’ according to a set of heuristics (instruction #13) and generates the triples describing the ontology.

Klink-2 was evaluated on a gold standard ontology including 88 research topics in the field of Semantic Web, which was manually generated by three senior researchers. It significantly outperformed the alternative algorithms ($p = 0.0005$), yielding a precision of 86% and a recall of 85.5%. More details about Klink-2 and its evaluation can be found in Osborne and Motta [37].

We selected Klink-2 among the other previously discussed solutions for a number of reasons. First, it is the only approach to our knowledge that was specifically designed to generate taxonomy of research areas. Secondly, it was already integrated and evaluated on a dump of the Scopus dataset, which we adopted in this study, yielding excellent performance on the fields of artificial intelligence and semantic web [37]. In third instance, it permits to define a number of pre-determinate relationships as basis for a new taxonomy. In particular, a human user can define a subsumption relation (i.e., *skos:broaderGeneric*), a *relatedEquivalent* one, or specify that two concepts should not be in any relationships. This functionality allows us to easily incorporate expert feedback in the ontology learning process. Therefore, the next iterations of the ontology will benefit from the knowledge of previous reviewers. We ran Klink-2 on the selected dataset, giving as initial seed the keyword “Software Engineering” and generated an OWL ontology of the field including 956 concepts and 5,461 relationships. We then selected the sub-branch of Software Architecture comprising 46 research areas and 71 terms (some research areas have multiple labels, such as “component based software” and “component-based software”).

4. Ontology refining. We generated a spreadsheet, containing the Software Architecture (SA) ontology as a tree diagram.²⁴ In this representation each concept of the ontology was illustrated by its level in

²⁴<http://tinyurl.com/yal6h3wu>

the taxonomy, its labels, and the number of papers annotated with the concepts. We also included a list of the 500 more popular terms in the papers that contained the keywords “Software Architecture” and “Software Engineering”, to assist the experts in remembering other concepts or terms that the algorithm may have missed.

We sent it to three senior researchers and asked them to correct the ontology as discussed in Section 4.1. The task took about 20 minutes and produced three revised spreadsheets. The feedback from the experts was integrated in the final ontology.²⁵ In case of disagreement we went with the majority vote.

The most frequent feedback regarded: (1) the deletion sub-areas that were incorrectly classified under SA (e.g., “software evolution”), (2) the introduction of sub-areas that were neglected by Klink-2 (e.g., “architecture concerns”), and (3) the inclusion of alternative labels for some category (e.g., alternative ways to spell “component-based architecture”).

5. Selection of primary studies. We then selected from the initial Scopus dump two datasets of primary studies to investigate the SA area: (1) **DSA** (Dataset SA, 3,467 publications), including all papers in the Scopus dataset that contain the terms “Software Architectures” or “Software Architecture” and include at least one of the subtopics of Software Architecture in the domain ontology, and (2) **DSA-MV** (Dataset SA – Main Venues, 1,586 publications), containing all the papers published in a list of well-known conferences and journals in the SE fields and in a particular in the SA area (see Table 1) and including at least one of the sub-topics of SA in the OWL ontology. We considered these two datasets since it may be interesting to analyze the discrepancy between generic SA papers and papers published in the main venues.

6. Classification of primary studies. We defined the mapping function as follows. A paper was classified under a certain category (e.g., service-oriented architectures) if it contained in the title, abstract

Table 1
List of venues used for the DSA-MV dataset

Conferences
WICSA - IEEE/IFIP Conference on Software Architecture, ECSA - European Conference on Software Architecture, CBSE - Int. ACM SigSoft Symposium on Component-based Software Engineering, QoSA - Conference on the Quality of Software Architecture, ICSE - ACM/IEEE Int. Conference on Software Engineering, ASE - IEEE/ACM Int. Conference on Automated Software Engineering, ESEC/FSE - European Software Engineering Conference, SEAA - Euromicro Conference on Software Engineering and Advanced Applications, ACM/SAC - ACM Symposium on Applied Computing
Journals
CACM - Communications of the ACM, ACM TOSEM - ACM Transactions on Software Engineering and Methodology, IEEE TSE - IEEE Trans. on Software Engineering, IEEE Software, Elsevier JSS - Journal of Systems and Software, Elsevier IST - Information and Software Technology, Wiley JSME/JSEP - Journal of software: Evolution and Process

²⁵<http://rexplore.kmi.open.ac.uk/data/edam/SE-ontology.owl>

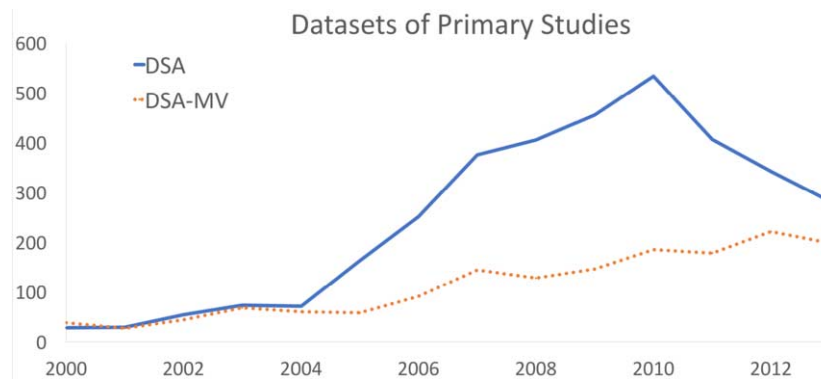


Fig. 4. Number of publications in DSA and DSA-MV over the years.

or keywords: (1) the label of the category (e.g., “service-oriented architectures”), (2) a *relevantEquivalent* of the category (e.g., “service oriented architecture”), (3) a *skos:broaderGeneric* of the category (e.g., “microservices”), or (4) a *relevantEquivalent* of any *skos:broaderGeneric* of the category (e.g., “microservice”). The advantage of this solution is that it allows us to map each category to a list of terms that can be automatically searched in the metadata of the papers. Therefore, the classification step can be handled automatically. In addition, it allows us to associate multiple categories to the same paper.

We chose this straightforward approach instead of other more complex methods based on word embeddings and string similarity [53], since it is simple to reproduce and yields the best precision, as discussed in Section 5.2. There we discuss also some recent approaches [53,55] yielding a more comprehensive set of topics and therefore a better recall, and illustrate how the choice of the method ultimately depends on the requested tradeoff between precision and recall.

In practice, we indexed titles, abstracts and keywords in an ElasticSearch²⁶ instance and we ran a PHP script that imported the ontology, performed the relevant queries on the metadata, and saved the result in a MariaSQL database.²⁷

7. Data synthesis. Fig. 4 shows the number of primary studies in the DSA and DSA-MV datasets. The DSA dataset follows the trend of the “Software Architecture” keyword in the Scopus dataset and decreases after 2010. Conversely, the size of DSA-MV grows steadily with the number of relevant conferences and journals.

We identified the main trends by running a script to count the number of studies about each subtopic in each year. Since the focus of the paper is the EDAM methodology, rather than a comprehensive analysis on these research sub-areas, we will briefly discuss only the main trends associated with the more popular subtopics (in terms of number of papers). The full results of this example study, however, are available at rexlore.kmi.open.ac.uk/data/edam and on Zenodo²⁸ and can be reused for supporting a more in-depth analysis of the field.

Figure 5 displays the number of publications and citations associated with the most popular sub-areas of SA. The papers in DSA yield on average 4.8 ± 2.1 in citations versus the 13.6 ± 7.0 citations of those in DSA-MV. Reasonably, this tendency suggests that the papers published in the main SA venues tend to be more recognized by the research community.

²⁶<https://www.elastic.co/>

²⁷<https://mariadb.org/>

²⁸<https://doi.org/10.5281/zenodo.2653924>

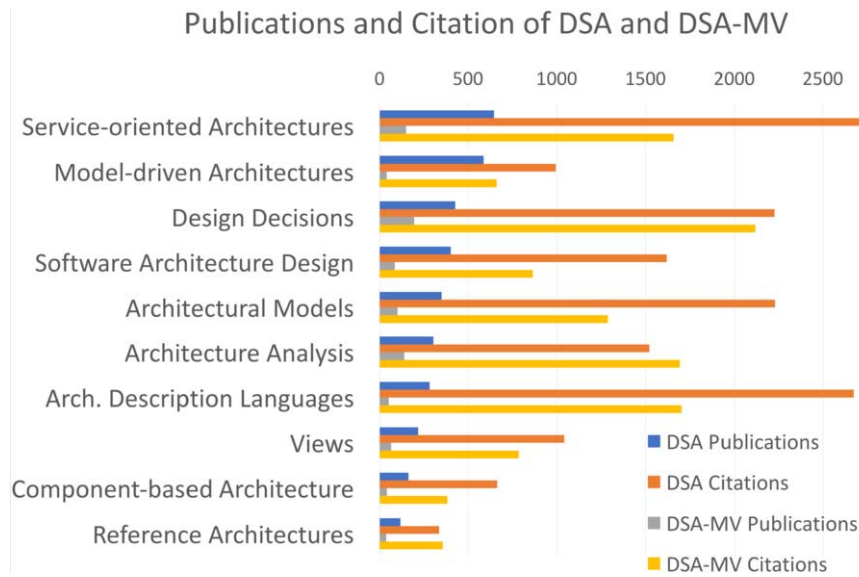


Fig. 5. Number of publications and citations of the main topics in DSA and DSA-MV.

Figure 6 shows the percentage of papers published over time in the main topics within SA. We focus on the 2005–2013 period, since in this interval the number of publications is high enough to highlight the topic trends.

Software-oriented Architectures appears to have been the most prominent topic before 2009, while from 2010, Model-driven Architectures appears to be the most popular topic in this dataset. We can also appreciate the rising of Design Decisions, that seems the most significant positive trend of the last period together with Architecture Description Languages.

Interestingly, the dataset regarding the main venues (DSA-MV) exhibits some different dynamics. Figure 7 highlights the difference between DSA and DSA-MV by showing for each topic the ratio between its number of publications and the total publications in the ten main topics. The research areas of Design Decisions and Views appear much more prominent in the main venues, while Model-Driven Architectures and Architecture Analysis are more popular in DSA. We can further analyze these differences by considering the main trends of the DSA-MV dataset, displayed by Fig. 8. The trend of Design Decisions in DSA-MV mirrors the one exhibited in DSA, both growing steadily from 2010. Conversely, Service-oriented Architectures, which has a negative trend in DSA, remains stable in DSA-MV.

5. Evaluation and discussion

In the following, we reflect on this preliminary application of EDAM. This section includes (1) a evaluation of our method versus six human annotators, (2) a comparison of several approaches for classifying primary studies, (3) an analysis of EDAM limitations, (4) a discussion about the implications for systematic mappings in Software Engineering, and (5) a discussion on how to reuse EDAM for other SRs.

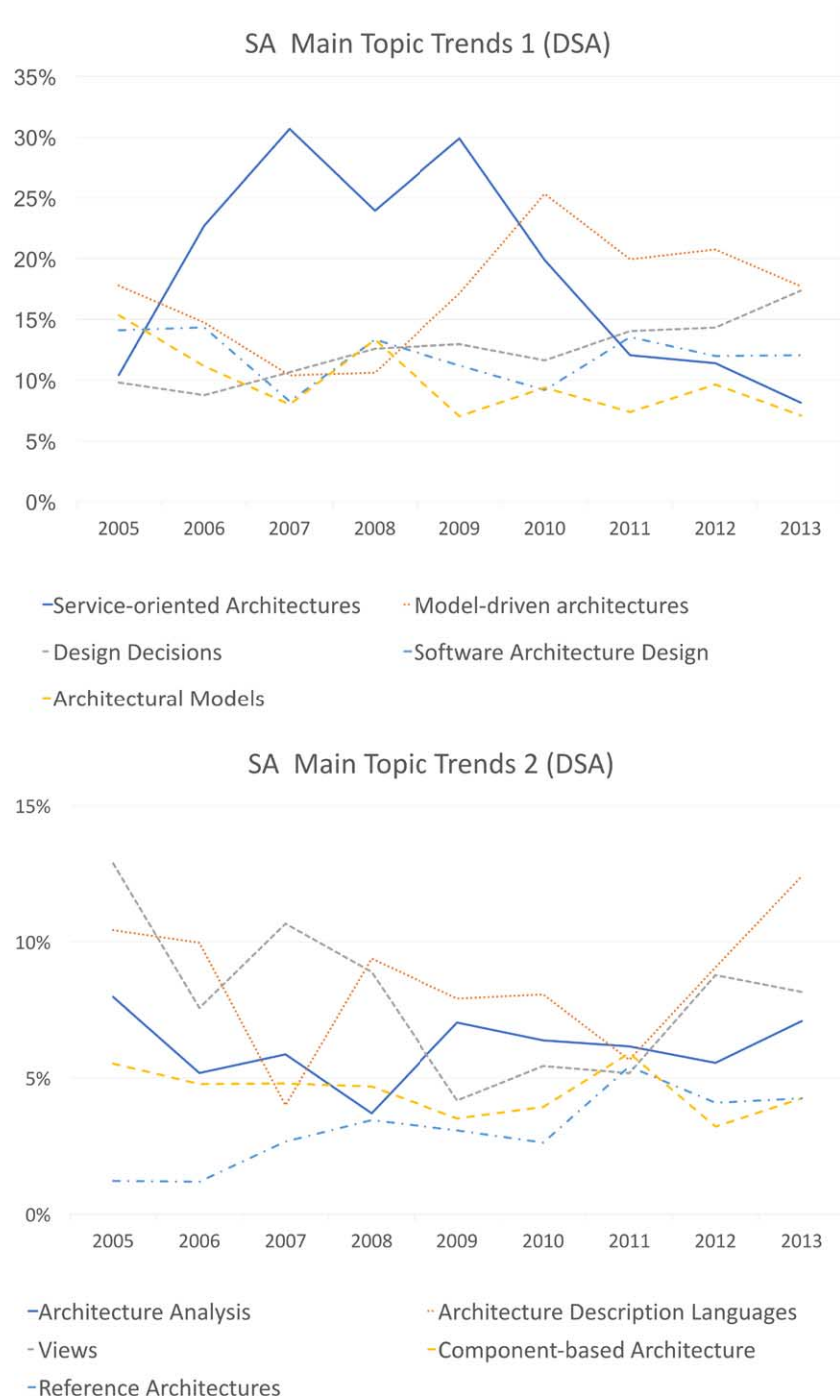


Fig. 6. Number of publications of the top ten main topics in DSA over time.

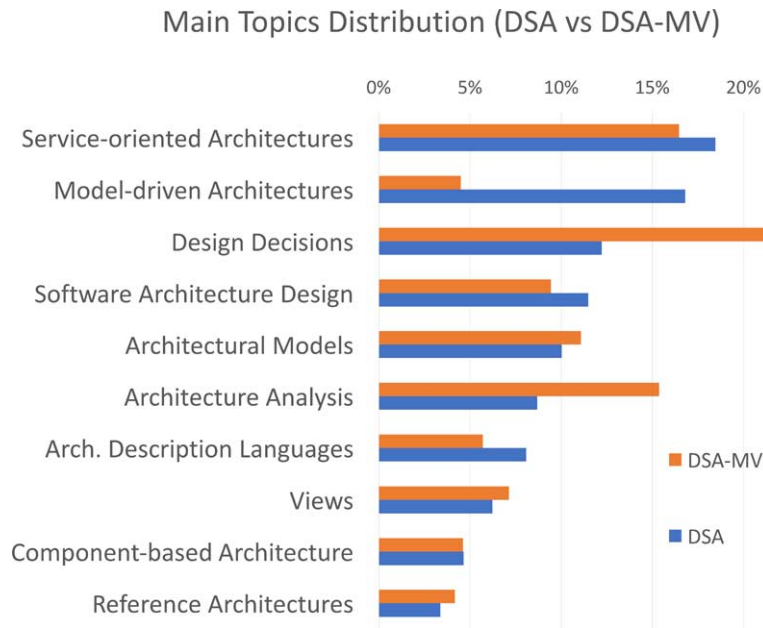


Fig. 7. Comparison DSA and DSA-MV in terms of topic distribution. The percentage value refers to the ratio between the number of publications in a topic and the total publications in the ten main topics.

5.1. Evaluation of the primary study classification

The most critical step of EDAM is the classification of primary studies. When these are correctly associated to the relevant topics, the subsequent analysis presents a realistic assessment of the landscape of the studied research field. Thus, even if working on a large number of papers can alleviate the weight of some minor misclassification mistakes, we need to be able to trust that the automatic classification process will obtain an accuracy similar to that yielded by human annotators.

We evaluated the ability of EDAM to correctly discriminate between different topics in the field of Software Architecture by (1) randomly selecting a set of 25 papers in the DSA dataset, (2) classifying them both with EDAM and with six human experts (researchers in the field of SA), and (3) comparing the results. For simplifying the task and allowing to compare the annotation algorithmically, we first selected five unambiguous categories from the main topics of SA: Design Decisions, Service-oriented Architectures, Model-driven Architectures, Architecture Description Languages, and Views. For each category, we randomly selected from the DSA dataset five primary studies that were classified by EDAM exclusively under that topic, for a total of 25 papers. These papers were described in a spreadsheet by means of their title, author list, abstract, and keywords. The human experts were given this spreadsheet and asked to classify each paper either with one of the five categories or with a “none of the above” tag. We then compared the seven annotation sets produced by the six human experts and by EDAM, considered as an additional annotator.²⁹

Table 2 shows the agreement between the annotators. It was computed by calculating the ratio of papers which were tagged with the same category by both annotators. EDAM has the highest average

²⁹The material and the results of the evaluation are available at <https://doi.org/10.5281/zenodo.2653924>.

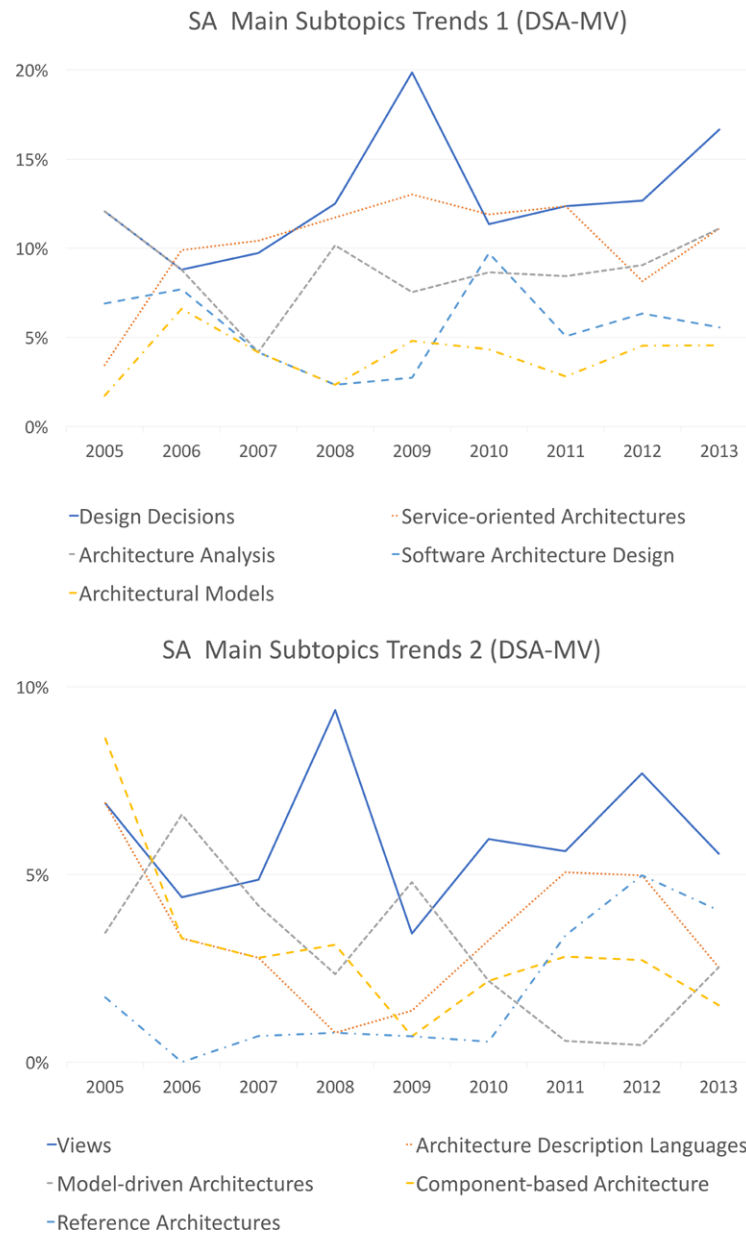


Fig. 8. Number of publications of the top ten main topics in DSA-MV over time.

agreement and it also yields the highest agreement with three out of six users. User5 does even better in this regards and has the highest agreement with four annotators.

The chi-square test run on the human users shows that their behaviours are significantly different ($p = 0.017$). However, if we group together users {2, 3, 5, 6} and users {1, 4}, we find no significant differences in the behaviour within each group ($p = 0.81$, $p = 0.38$, good intra-group agreement), while there are between the two groups ($p = 0.0007$).

Interestingly, users {1, 4} were two students at the beginning of their PhD, hence still relatively new to

Table 2

Agreement between annotators (including EDAM) and average agreement of each annotator. In bold the best agreements for each annotator

	EDAM	User1	User2	User3	User4	User5	User6
EDAM		56%	68%	64%	64%	76%	64%
User1	56%		40%	56%	36%	48%	44%
User2	68%	40%		64%	52%	76%	64%
User3	64%	56%	64%		52%	64%	68%
User4	64%	36%	52%	52%		64%	52%
User5	76%	48%	76%	64%	64%		72%
User6	64%	44%	64%	68%	52%	72%	
Av. Agreement	66%	45%	58%	59%	51%	63%	60%

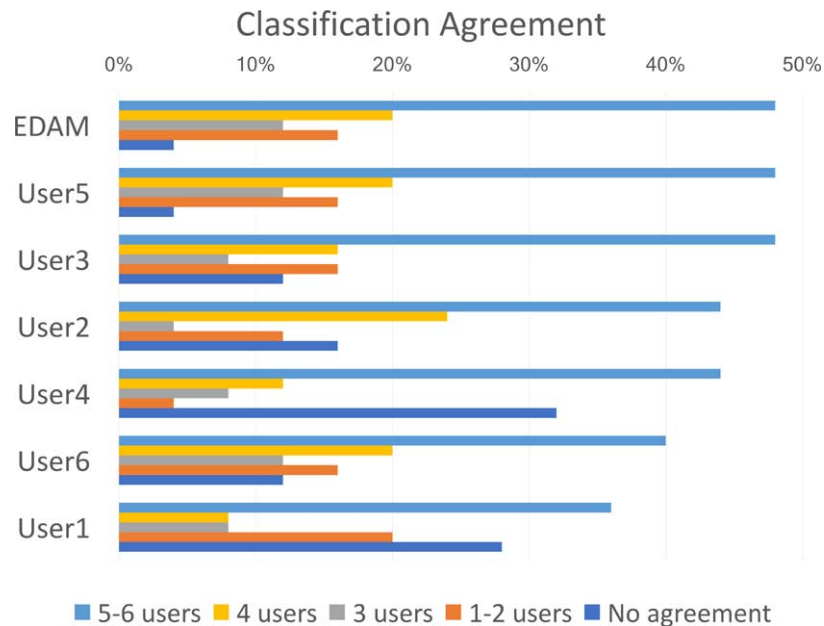


Fig. 9. Percentage of annotations that agree with other n annotators.

the domain. This could suggest the importance of considerable domain experience for this task. EDAM exhibits a behaviour consistent with the most senior group, from which it is not significantly different ($p = 0.77$).

As anticipated, a good way to measure the performance of annotators is their agreement with the majority of other expert users. Figure 9 shows the percentage of annotations of each annotator that agree with other n annotators. EDAM agrees with four out of six human annotators for 68% of the studies, it agrees with at least three of them for 80% of the studies, and it agrees with at least one of them for all the studies but one. Indeed, the categories generated by EDAM coincide with the ones suggested by the relative majority of users in 84% of the cases. Therefore, EDAM's performance is comparable to the performance of the two annotators (User5 and User3) with the highest agreement with the user majority.

We further confirmed these findings by computing the Cohen's kappa between each couple of annotators and between each of the annotators and EDAM. The inter-annotator agreement was 0.57, typically indicating a moderate agreement [27]. The average agreement of EDAM with the annotators was 0.58, confirming that this method performs in line with the annotators. In addition, we note that EDAM always agrees with the majority for the studies in which no more than one annotator disagrees. It thus seems to perform well in handling simple not-ambiguous papers, that nonetheless human experts may sometimes get wrong.

In conclusion, this study suggests that the EDAM classification step generates annotations that agree with the majority of human experts and are not statistically different from the ones produced by the senior group.

Naturally, EDAM performance may change according to the quality of the ontology and the domain knowledge of the human users that refined it. EDAM is not an alternative to human experts, rather a methodology that allows humans to annotate on a larger scale, by defining a sound domain knowledge and a mapping function. However, this preliminary example application already shows very promising results.

5.2. *Comparison of classifiers for primary studies*

EDAM can adopt many different approaches for automatically classifying primary study. The choice of the method ultimately depends on its effectiveness of the available approaches on the domain under analysis and on the preferred precision-recall tradeoff. For instance, in a prevalently automatic mapping study on a large set of papers, precision is of paramount importance, and the missed topics may be compensated by the numerosity of the sample. Conversely, when the results are validated in some way by human experts (e.g., [39]) or when the goal is to detect emerging trends that may appear in few publications, a better strategy may be to sacrifice some precision for producing a more comprehensive set of topics.

In this section we compare several approaches that can be used with EDAM to automatically classify studies according to a taxonomy of research topics, and discuss their tradeoff. We focus on unsupervised approaches, since typically the authors of a systematic review do not have the resources to prepare a large gold standard to train a supervised classifier on a potentially new taxonomy.

We evaluated seven alternative approaches on a gold standard of 70 papers [53] within the fields of Semantic Web (23 papers), Natural Language Processing (23), and Data Mining (24). These papers were selected by retrieving the most cited papers from Microsoft Academic Graph containing in the title or the abstract those relevant fields. Each paper was annotated by three domain experts (for a total of 21 different annotators) and was associated with 14.4 ± 7.0 topics using majority vote in case of disagreement. The inter-annotator agreement was 0.45 ± 0.18 according to Fleiss' Kappa, which indicates a moderate inter-rater agreement [27]. The data produced in the evaluation and the Python implementation of the approaches are available at <https://cso.kmi.open.ac.uk/cso-classifier/>. A more comprehensive version of this analysis, with additional baselines that are out of scope for this paper (e.g., not producing a set of pre-defined categories), is available in Salatino et al. [53].

All the tested classifiers analysed the title and abstract of the 70 papers and assigned them with a set of topics drawn from the Computer Science Ontology (CSO), a recently released taxonomy of research areas [54]. This knowledge base covers well the three mentioned fields, including a total of 35 sub-topics for the Semantic Web, 173 for Natural Language Processing, and 396 for Data Mining. LDA100, LDA500, and LDA1000 are based on a Latent Dirichlet Allocation (LDA) model [8] trained over 4.6M

papers in Computer Science from Microsoft Academic Search. Specifically, LDA100 used a model trained with 100 topics, LDA500 on 500 topics, and LDA1000 on 1,000 topics. Similarly to [7], these classifiers generate a set of CSO topics from the LDA topics by first producing a set of topics with a probability of at least j and all their terms with a probability of at least k . Then they map these terms to CSO by returning all CSO topics having Levenshtein similarity higher than 0.8 with them. We performed a grid search for finding the best values of j and k on the gold standard and report here the best results of each classifier in term of F-measure.

TF-IDF produces a ranked list of terms using TF-IDF [44] (the IDF was computed on the same set adopted for LDA) and all the CSO topics having Levenshtein similarity higher than 0.8 with the first 30 terms.

Direct Mapping (DM) is the approach used for the implementation of EDAM described in step 6 of Section 4.2. This same method was also used by the first version of the Smart Topic Miner (STM) [39], the system adopted by Springer Nature to classifying proceedings in the field of Computer Science. It returns all topics that explicitly appear in the papers or that are entailed by the ones appearing in the paper according to the ontology.

The CSO Classifier v.1 (CSO-C1) is an unsupervised approach presented in Salatino et al. [55] that extracts a combination of n-grams (unigrams, bigrams, and trigrams) from the text and returns all the topics that have a Levenshtein similarity higher than t ($t = 0.94$ as in the implementation reported in Salatino et al. [55]). Finally, the CSO Classifier v.2 (CSO-C2) [53] is a recent evolution of the previous classifier, which uses part-of-speech tagging to identify promising terms and then exploits word embeddings to infer semantically related topics that may not explicitly appear in the paper. This solution has also been adopted by the current version of the Smart Topic Miner [51].

Table 3 reports on the resulting values of precision, recall and F-measure. The approaches based on LDA performed quite poorly. An analysis of the results revealed that the topics returned by the models are both noisy and coarse-grained, often clustering together distinct topics from CSO. Indeed, while LDA works quite well at identifying the main topics of a large collection of documents, it does not traditionally perform equally well when characterizing specific research topics, which may be associated with a relatively low number of publications (50–200), as discussed in [36]. The approaches based on TF-IDF worked slightly better, yielding an F-measure of 30.1%. DM, used in our exemplary EDAM implementation, yielded the best precision of all the approaches (80.8%). Indeed, this method focuses on topics that are explicitly mentioned in the text, which tends to be very relevant. However, this solution

Table 3

Values of precision, recall, and F-measure for the seven classifiers. In bold the best results

Approach	Description	Precision	Recall	F-Measure
LDA100	LDA [6], 100 topics, 4.6M papers, Lev. 0.8	9.2%	17.4%	12.1%
LDA500	LDA [6], 500 topics, 4.6M papers, Lev. 0.8	9.1%	23.4%	13.1%
LDA1000	LDA [6], 1000 topics, 4.6M papers, Lev. 0.8	12.0%	11.5%	11.7%
TF-IDF	TF-IDF [36], 4.6M papers, Lev. 0.8	40.4%	24.0%	30.1%
DM	Direct Match [32], described in Section 4.2 and used for the EDAM implementation in this paper	80.8%	58.2%	67.6%
CSO-C1	The CSO classifier ver. 1 [44]	78.3%	63.8%	70.3%
CSO-C2	The CSO classifier ver. 2 [45]	73.0%	75.3%	74.1%

naturally obtains a relatively low recall of 58.2%. CSO-C1, which expands the set of terms by considering string similarity, obtained a better recall (63.8%) but a lower precision (78.3%). Finally, CSO-C2 yielded the best performance in term of both recall (75.3%) and F-measure (74.1%).

DM performed significantly better ($p < 10^{-7}$ with the McNemar's test) than the first four approaches. In turn, CSO-C2 performed significantly better ($p < 10^{-7}$) than DM and CSO-C1.

5.3. Limitations

In this section we discuss EDAM limitations based on the categorization given in Wohlin et al. [66].

For *internal validity* we have identified two main threats that regard the generation of a reliable ontology, which is key to select relevant studies that directly fulfill the selection criteria (and hence correspond to the primary studies for the study at hand). In particular:

Ontology learning (step 3): hierarchy is important. The domain ontology, automatically inferred by the ontology learning technique, is structured hierarchically. Therefore, an area marked as *subarea* (e.g., architecture description languages) is subsumed by the previous area at the upper level of the taxonomy (e.g., Software Architecture). *Deeper hierarchies bring finer-grained topics, and therefore a higher precision in the classification process.*

During the application of ontology learning techniques to various research areas (not reported in this paper for the sake of brevity) we found that current ontology learning methods usually identify only mature (in terms of number of publications) research areas. Emerging topics may be excluded, thus reducing the granularity of recent fields' ontologies.

To alleviate this problem, human experts may be asked to manually identify the most recent areas and to possibly adopt ontology forecasting techniques [9]. Therefore, the role of experts in improving the quality and deepness of the hierarchy is indeed critical. For the sake of this study, aimed at showing the advantages of automation, the relatively small number of experts was acceptable. However, a larger and more diversified pool of experts should be involved when the research area under investigation is broader.

Ontology refinement (step 4): experience matters. As illustrated in Fig. 2, EDAM requires human expertise to refine the automatically generated ontology (step 4). This task is not always straightforward, since humans can have different views on the foundational conceptual elements characterizing a certain discipline. Those differences may be related to many factors, such as the researcher's exposure to the research area under investigation, seniority, broad vs. specialized knowledge on specific sub-disciplines. Our preliminary experiments allow us to conclude that senior domain experts, with a mature yet wide view on the research area under investigation, should be selected to minimize this threat.

The main threats for *external validity* regard the practical exploitation of EDAM. In particular:

Scholarly dataset: different research areas require different datasets. This paper reports on our experience with EDAM's application to the Software Architecture research area. Since the domain of Software Engineering is well represented in the Scopus dataset, we are not facing generalizability issues. However, moving to a totally different domain would require taking into account (assuming to have access to) different scholarly datasets.

Unfortunately, finding up-to-date datasets of scholarly data covering the field under analysis is not always easy and this could be a threat to our approach. Nonetheless, the movement toward

open access is helping in mitigating this issue by making available a variety of datasets containing machine-readable data about scientific publications, e.g., Microsoft Academic Graph³⁰ [58], CORE³¹ [24], OpenCitations³² [40], DBLP³³ [28], Bio2RDF³⁴ [6], ScholarlyData.org³⁵ [19], Nanopub.org³⁶ [25], Semantic Scholar,³⁷ and others.

Tool support: closed-source tools. EDAM is making use of some closed-source, proprietary tools for running some of the tasks. This may reduce the application of our approach from other research groups. In order to mitigate this threat, we are planning to release a web service accessible by other colleagues interested to carry out an EDAM study.

Research Questions: some may not be automatized. Many research questions that are typical of mapping studies can be answered by producing relevant analytics [66], e.g., by counting the number of publications, authors, and venues associated with certain topics in subsequent years. However, some more complex research questions may still require domain experts to manually analyse the relevant studies, e.g., for classifying them in categories that a state of the art classifier would be unable to detect with good accuracy. This is an inherent limitation of the methodology. Nonetheless in many of these cases a preliminary classification by an automatic system may still alleviate the expert work load, e.g., by reducing the set of publications that need to be manually analysed. In addition, the performance of entity extraction and linking tools is steadily improving [4,18,47], allowing to extract increasingly better representations of research knowledge from scientific articles. Therefore, the number of research questions that can be addressed algorithmically may increase over the following years.

5.4. Implications for systematic mappings

There are a few implications that can potentially change the way we perform systematic mapping studies in Software Engineering. As mentioned in Section 4, these implications regard:

Scalability: size does not matter anymore. EDAM can process a potentially endless set of publications. This allows e.g., mapping studies to be based on *all* relevant primary studies, previously scoped down due to the fact that humans could not manually process hundreds or thousands of papers.

Objectivity: the automatic classification is less biased. The automatic classification of primary studies does not suffer from the biases of specific human annotators. Nonetheless, the quality of the classification appears on par with the one produced by the human annotators.

Reproducibility: study duplication and extension is easy. Thanks to EDAM, replicating or extending studies, either by the same researcher or by someone else, requires simple tuning, e.g., to extend the publication period, or to select different views illustrating the publication trends of interest.

³⁰<https://academic.microsoft.com/>

³¹<https://core.ac.uk>

³²<http://opencitations.net/>

³³<http://dblp.uni-trier.de/>

³⁴<http://bio2rdf.org/>

³⁵<http://www.scholarlydata.org/>

³⁶<http://nanopub.org/>

³⁷<https://www.semanticscholar.org/>

Granularity of the study: zooming-in and -out is simpler. Thanks to the fact that the selection and classification of primary studies is based on an domain ontology, and of course to automation, EDAM allows to tune the depth of the classification the researcher desires in a given research area. Such tuning just requires setting the level of categories and sub-categories to be included in the classification, and then re-run the methodology.

5.5. Reusing EDAM for other systematic reviews

EDAM can be applied to any domain of interest and for different types of studies. The scenarios that we envisage are discussed below and illustrated in Fig. 10. They are: S1) Application of EDAM to a *new* application domain, S2) Mapping study *replication*, S3) Mapping study *refinement*, and S4) *Systematic literature review*.

Application of EDAM to a new application domain (S1). In the basic scenario (S1), the ontology for the new application domain is not yet available. In this case, the complete process illustrated in Fig. 2 (and emphasized in Fig. 10(S1)) shall be applied. This is the scenario followed in the work presented in this article. It is applicable while investigating a new domain notwithstanding its specific characteristics.

If instead a researcher wants to perform a SR in a domain for which the ontology already exists (scenario S2), such generated domain ontology can be *reused* in the following two ways, depending on the specific study goal:

Mapping Study Replication (same classification, S2a). Suppose we want to replicate a pre-existing EDAM mapping study conducted at time t_0 , in order to update the list of primary studies and related analysis at time t_1 (e.g., update in year 2020 the study on Software Architecture presented in this paper). In this case, we can directly reuse the previously generated ontology (cf. Fig. 10(S2a)). The list of (updated) primary studies can be automatically re-calculated (in step 5) and used (in step 6) for classification and analysis purposes. Notice, however, that this scenario does not address the potential need to *update* the list of topics. Such a scenario is covered below.

Mapping Study Replication (updated classification, S2b). Differently from scenario S2a, we may be interested to replicate a pre-existing study *and* also include any new topics that may have emerged in the period between time t_0 and time t_1 (e.g., updating this study in year 2020 while including new topics appeared after this study). This need requires an update of the domain ontology; therefore, the process in Fig. 10(S2b) must be run from step 4 onward.

Another scenario (S3) accommodates the case in which we want to *refine* the classification and analysis conducted as a mapping study. In the current approach, as shown in the Software Architecture domain scenario, step 5 in Fig. 2 returns a set of primary studies that can be further classified into sub-domains (e.g., Architectural Styles, being one element of our ontology, can be further refined to discover all the papers that cover selected styles). We identify two sub-scenarios in order to provide a refinement of sub-domains contents:

Mapping Study Refinement with classic selection criteria (S3a). In this scenario, one may classify the articles into sub-domains of interest by applying the inclusion and exclusion criteria [23] to the primary studies selected in step 5 of EDAM. For example, knowing that Publish-Subscribe, Client-Server, and Event-driven are sub-domains of Architectural Styles, we introduce selection

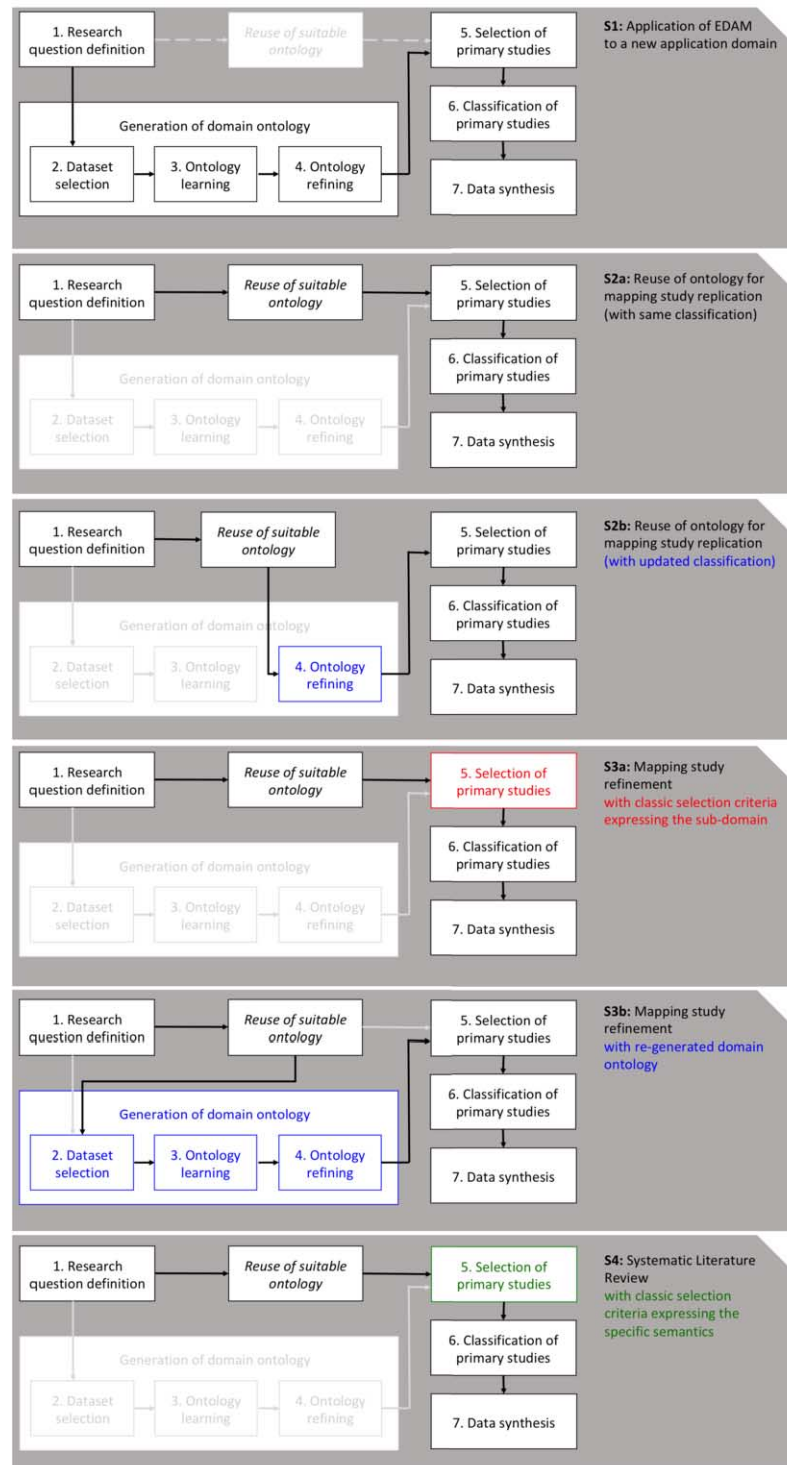


Fig. 10. Possible EDAM applications.

criteria to position Architectural Styles articles into those categories. This approach allows us to zoom into a specific sub-domain of interest and extract the articles fitting in the specific target sub-domain.

Mapping Study Refinement with re-generated domain ontology (S3b). The selected sub-domain of interest may contain hundreds of papers (for example, the Design Decisions sub-domain in our study includes 428 papers). Consequently, applying the selection criteria reported in scenario S3a may be cumbersome, requiring the manual analysis of most of those papers. Alternatively, the researcher may execute an additional round of steps 2–4 to refine the domain ontology for the specific sub-domain (cf. Fig. 10(S3b)). This scenario is similar to S1, but applied to a specific sub-domain of interest.

A fourth scenario sees the researcher is interested to run a systematic literature review (SLR) on specific research questions:

Systematic Literature Reviews (S4). In step 5 (cf. Fig. 10(S4)), given the list of primary studies generated based on the existing ontology, we may run the *classic* SLR approach [22] to select those papers that fit with the research questions of interest. Differently from scenario S3a, S4 adds the semantics beyond the definition of the domain, and encapsulated into the research questions and the corresponding selection criteria. E.g., given the list of all studies on Software Architecture styles, one may want to perform an SLR to analyze those approaches that are adopted in industrial settings.

6. Conclusions and future work

In this paper we have presented EDAM, an expert-driven automated methodology to assist systematic reviews. Its application to the Software Architecture research area shows preliminary and very promising results.

Motivated by the large amount of time and effort needed by classic methodologies to select and classify the primary studies, EDAM offers benefits that can help SE researchers to dedicate most of their time to the most cognitive-intensive tasks like e.g., interpretation of the trends and extraction of lessons and research gaps.

Additional benefits have been emphasized in Section 4.1 (after presenting EDAM) and Section 5.4 (discussing implications for systematic mappings). Among the benefits we mention the great potential for re-using EDAM and in particular domain ontologies and functions to build a shared framework helping the research community at large. Much can be done in this direction.

Our next step is to complement EDAM with automated forward snowballing to further reduce the effort for identifying relevant primary studies. With the same goal, we are planning to investigate other possible data synthesis techniques through machine learning techniques or the (manual) intervention of human experts. Last, but most important for us, we plan to reconstruct the 25 years of the Software Architecture body of knowledge by fully exploiting EDAM automation and human expertise.

Acknowledgements

The authors would like to thank the colleagues which donated their time and expertise by contributing to this study as domain experts and/or annotators: Paris Avgeriou, Barbora Buhnova, Rafael Capilla, Jan

Carlson, Ivica Crnkovic, John Grundy, Rich Hilliard, Heiko Kozirolek, Anton Jansen, Ivano Malavolta, Leonardo Mariani, Marina Mongiello, Matthias Naab, Patrizio Pelliccione, Mohammad Sharaf, Damian Andrew Tamburri, Antony Tang, Jan Martijn van der Werf, Smrithi Rekha Venkatasubramanian, Rainer Weinreich, Danny Weyns, Eoin Woods, and Uwe Zdun.

We also thank Davide Falessi for reviewing an earlier version of this manuscript, and Elsevier BV for providing us with access to its large repository of scholarly data.

References

- [1] A. Al-Zubidy, J.C. Carver, D.P. Hale and E.E. Hassler, Vision for SLR tooling infrastructure: Prioritizing value-added requirements, *Information and Software Technology* **91** (2017), 72–81. doi:[10.1016/j.infsof.2017.06.007](https://doi.org/10.1016/j.infsof.2017.06.007).
- [2] R. Alghamdi and K. Alfalqi, A survey of topic modeling in text mining. I, *J. ACSA* **6**(1) (2015), 147–153. doi:[10.14569/IJACSA.2015.060121](https://doi.org/10.14569/IJACSA.2015.060121).
- [3] E. Aromataris and D. Riitano, Constructing a search strategy and searching for evidence, *American Journal of Nursing* **114**(5) (2014), 49–56. doi:[10.1097/01.NAJ.0000446779.99522.f6](https://doi.org/10.1097/01.NAJ.0000446779.99522.f6).
- [4] I. Augenstein, S. Padó and S.R. Rudolph, Lodifier: Generating linked data from unstructured text, in: *Extended Semantic Web Conference*, Springer, 2012, pp. 210–224. doi:[10.1007/978-3-642-30284-8_21](https://doi.org/10.1007/978-3-642-30284-8_21).
- [5] F.B.V. Benitti, Exploring the educational potential of robotics in schools: A systematic review, *Computers & Education* **58**(3) (2012), 978–988. doi:[10.1016/j.compedu.2011.10.006](https://doi.org/10.1016/j.compedu.2011.10.006).
- [6] F. Belleau, M.-A. Nolin, N. Tourigny, P. Rigault and J. Morissette, Bio2rdf: Towards a mashup to build bioinformatics knowledge systems, *Journal of biomedical informatics* **41**(5) (2008), 706–716. doi:[10.1016/j.jbi.2008.03.004](https://doi.org/10.1016/j.jbi.2008.03.004).
- [7] S. Bhatia, J.H. Lau and T. Baldwin, Automatic labelling of topics with neural embeddings, 2016, arXiv preprint [arXiv:1612.05340](https://arxiv.org/abs/1612.05340).
- [8] D.M. Blei, A.Y. Ng and M.I. Jordan, Latent Dirichlet allocation, *Journal of machine Learning research* **3** (2003), 993–1022. <https://dl.acm.org/citation.cfm?id=944937>.
- [9] A.E. Cano-Basave, F. Osborne and A.A. Salatino, Ontology forecasting in scientific literature: Semantic concepts prediction based on innovation-adoption priors, in: *Knowledge Engineering and Knowledge Management: 20th International Conference, EKAW 2016, Proceedings 20*, Bologna, Italy, November 19–23, 2016. Springer, 2016, pp. 51–67. doi:[10.1007/978-3-319-49004-5_4](https://doi.org/10.1007/978-3-319-49004-5_4).
- [10] M.J.M. Chinapaw, K.I. Proper, J. Brug, W. Van Mechelen and A.S. Singh, Relationship between young peoples' sedentary behaviour and biomedical health indicators: A systematic review of prospective studies, *Obesity reviews* **12**(7) (2011), e621–e632. doi:[10.1111/j.1467-789X.2011.00865.x](https://doi.org/10.1111/j.1467-789X.2011.00865.x).
- [11] P. Cimiano and J. Völker, *Text2onto*, Springer, 2005, pp. 227–238. doi:[10.1007/11428817_21](https://doi.org/10.1007/11428817_21).
- [12] F.Q.B. da Silva, M. Suassuna, A. César, C. França, A.M. Grubb, T.B. Gouveia, C.V.F. Monteiro and I.E. dos Santos, Replication of empirical studies in software engineering research: A systematic mapping study, *Empirical Software Engineer* **19**(3) (2014), 501–557. doi:[10.1007/s10664-012-9227-7](https://doi.org/10.1007/s10664-012-9227-7).
- [13] J.C. de Almeida Biolchini, P.G. Mian, A.C. Cruz Natali, T.U. Conte and G.H. Travassos, Scientific research ontology to support systematic review in software engineering, *Advanced Engineering Informatics* **21**(2) (2007), 133–151. doi:[10.1016/j.aei.2006.11.006](https://doi.org/10.1016/j.aei.2006.11.006).
- [14] R.M. De Mello and G.H. Travassos, Surveys in software engineering: Identifying representative samples, in: *Proceedings of the 10th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement, ESEM '16*, ACM, New York, NY, USA, 2016, pp. 55:1–55:6. doi:[10.1145/2961111.2962632](https://doi.org/10.1145/2961111.2962632).
- [15] T. Di Noia, R. Mirizzi, V.C. Ostuni, D. Romito and M. Zanker, Linked open data to support content-based recommender systems, in: *Proceedings of the 8th International Conference on Semantic Systems*, ACM, 2012, pp. 1–8. doi:[10.1145/2362499.2362501](https://doi.org/10.1145/2362499.2362501).
- [16] K.R. Felizardo, E. Mendes, M. Kalinowski, É.F. Souza and N.L. Vijaykumar, Using forward snowballing to update systematic reviews in software engineering, in: *Proceedings of the 10th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement*, ACM, 2016, p. 53. doi:[10.1145/2961111.2962630](https://doi.org/10.1145/2961111.2962630).
- [17] B. Filipovic, B. Van Lindschoten, G. Procaccianti and P. Lago, Systematic Literature Study on Sustainable Software, VU Technical Report, 2, 2017. https://research.vu.nl/ws/portalfiles/portal/58444907/Technical_Report_SLR.pdf.
- [18] A. Gangemi, V. Presutti, D. Reforgiato Recupero, A.G. Nuzzolese, F. Draicchio and M. Mongiovì, Semantic web machine reading with fred, *Semantic Web* **8**(6) (2017), 873–893. doi:[10.3233/SW-160240](https://doi.org/10.3233/SW-160240).
- [19] A.G. Nuzzolese, A.L. Gentile, V. Presutti and A. Gangemi, Conference linked data: The scholarlydata project, in: *International Semantic Web Conference*, Springer, 2016, pp. 150–158. doi:[10.1007/978-3-319-46547-0_16](https://doi.org/10.1007/978-3-319-46547-0_16).

- [20] E. Hassler, J.C. Carver, D. Hale and A. Al-Zubidy, Identification of SLR tool needs – results of a community workshop, *Information and Software Technology* **70** (2016), 122–129. doi:[10.1016/j.infsof.2015.10.011](https://doi.org/10.1016/j.infsof.2015.10.011).
- [21] B. Kitchenham, Procedures for performing systematic reviews, *Keele, UK, Keele University* **33**(2004) (2004), 1–26. <http://www.inf.ufsc.br/~aldo.vw/kitchenham.pdf>.
- [22] B. Kitchenham and P. Brereton, A systematic review of systematic review process research in software engineering, *Information and software technology* **55**(12) (2013), 2049–2075. doi:[10.1016/j.infsof.2013.07.010](https://doi.org/10.1016/j.infsof.2013.07.010).
- [23] B.A. Kitchenham and S. Charters, Guidelines for performing systematic literature reviews in software engineering, 2007. https://edisciplinas.usp.br/pluginfile.php/4108896/mod_resource/content/2/slrPCS5012_highlighted.pdf.
- [24] P. Knott and Z. Zdrahal, Core: Three access levels to underpin open access, *D-Lib Magazine* **18**(11/12) (2012). doi:[10.1045/november2012-knott](https://doi.org/10.1045/november2012-knott).
- [25] T. Kuhn, P.E. Barbano, M.L. Nagy and M. Krauthammer, Broadening the scope of nanopublications, in: *Extended Semantic Web Conference*, Springer, 2013, pp. 487–501. doi:[10.1007/978-3-642-38288-8_33](https://doi.org/10.1007/978-3-642-38288-8_33).
- [26] M. Kuhrmann, D. Méndez Fernández and M. Daneva, On the pragmatic design of literature studies in software engineering: An experience-based guideline, *Empirical Software Engineer* **22** (2017), 1–40. doi:[10.1007/s10664-016-9495-8](https://doi.org/10.1007/s10664-016-9495-8).
- [27] J.R. Landis and G.G. Koch, The measurement of observer agreement for categorical data, *Biometrics* **33**(1) (1977), 159–174. doi:[10.2307/2529310](https://doi.org/10.2307/2529310).
- [28] M. Ley, Dbp: Some lessons learned, *Proceedings of the VLDB Endowment* **2**(2) (2009), 1493–1500. doi:[10.14778/1687553.1687577](https://doi.org/10.14778/1687553.1687577).
- [29] X. Liu, Y. Song, S. Liu and H. Wang, Automatic taxonomy construction from keywords, in: *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ACM, 2012, pp. 1433–1441. doi:[10.1145/2339530.2339754](https://doi.org/10.1145/2339530.2339754).
- [30] A. Mannocci, F. Osborne and E. Motta, The evolution of IJHCS and CHI: A quantitative analysis, *International Journal of Human-Computer Studies* (2019). [arXiv:1908.04088](https://arxiv.org/abs/1908.04088).
- [31] C. Marshall, P. Brereton and B. Kitchenham, Tools to support systematic reviews in software engineering: A cross-domain survey using semi-structured interviews, in: *Proceedings of the 19th International Conference on Evaluation and Assessment in Software Engineering*, ACM, 2015, p. 26. doi:[10.1145/2745802.2745827](https://doi.org/10.1145/2745802.2745827).
- [32] P.N. Mendes, M. Jakob, A. García-Silva and C. Bizer, Dbpedia spotlight: Shedding light on the web of documents, in: *Proceedings of the 7th International Conference on Semantic Systems*, ACM, 2011, pp. 1–8. doi:[10.1145/2063518.2063519](https://doi.org/10.1145/2063518.2063519).
- [33] J.S. Moller, K. Petersen and E. Mendes, Survey guidelines in software engineering: An annotated review, in: *Proceedings of the 10th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement – ESEM '16*, ACM Press, 2016, pp. 1–6. doi:[10.1145/2961111.2962619](https://doi.org/10.1145/2961111.2962619).
- [34] E. Mourão, M. Kalinowski, L. Murta, E. Mendes and C. Wohlin, Investigating the use of a hybrid search strategy for systematic reviews, in: *Proceedings of the 11th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement, ESEM '17*, IEEE Press, 2017, pp. 193–198. doi:[10.1109/ESEM.2017.30](https://doi.org/10.1109/ESEM.2017.30).
- [35] R.F. Octaviano, K.R. Felizardo, J.C. Maldonado and C.P. Sandra, Semi-automatic selection of primary studies in systematic literature reviews: Is it reasonable?, *Empirical Software Engineer* **20**(6) (2015), 1898–1917. doi:[10.1007/s10664-014-9342-8](https://doi.org/10.1007/s10664-014-9342-8).
- [36] F. Osborne and E. Motta, Mining semantic relations between research areas, in: *International Semantic Web Conference 2012*, Springer, 2012, pp. 410–426. doi:[10.1007/978-3-642-35176-1_26](https://doi.org/10.1007/978-3-642-35176-1_26).
- [37] F. Osborne and E. Motta, Klink-2: Integrating multiple web sources to generate semantic topic networks, in: *International Semantic Web Conference 2015*, Springer, 2015, pp. 408–424. doi:[10.1007/978-3-319-25007-6_24](https://doi.org/10.1007/978-3-319-25007-6_24).
- [38] F. Osborne, E. Motta and P. Mulholland, Exploring scholarly data with rexplore, in: *International Semantic Web Conference 2013*, Springer, 2013, pp. 460–477. doi:[10.1007/978-3-642-41335-3_29](https://doi.org/10.1007/978-3-642-41335-3_29).
- [39] F. Osborne, A. Salatino, A. Birukou and E. Motta, Automatic classification of springer nature proceedings with smart topic miner, in: *International Semantic Web Conference 2016*, Springer, 2016, pp. 383–399. doi:[10.1007/978-3-319-46547-0_33](https://doi.org/10.1007/978-3-319-46547-0_33).
- [40] S. Peroni, A. Dutton, T. Gray and D. Shotton, Setting our bibliographic references free: Towards open citation data, *Journal of Documentation* **71**(2) (2015), 253–277. doi:[10.1108/JD-12-2013-0166](https://doi.org/10.1108/JD-12-2013-0166).
- [41] K. Petersen, R. Feldt, S. Mujtaba and M. Mattsson, Systematic mapping studies in software engineering, in: *Proceedings of the 12th International Conference on Evaluation and Assessment in Software Engineering, EASE*, British Computer Society, Swinton, UK, UK, 2008, pp. 68–77. <https://dl.acm.org/citation.cfm?id=2227115.2227123>.
- [42] K. Petersen, S. Vakkalanka and L. Kuzniarz, Guidelines for conducting systematic mapping studies in software engineering: An update, *Information and Software Technology* **64** (2015), 1–18. doi:[10.1016/j.infsof.2015.03.007](https://doi.org/10.1016/j.infsof.2015.03.007).
- [43] G. Petrucci, C. Ghidini and M. Rospocher, Ontology learning in the deep, in: *Knowledge Engineering and Knowledge Management: 20th International Conference, EKAW 2016, Proceedings 20*, Bologna, Italy, November 19–23, 2016, Springer, 2016, pp. 480–495. doi:[10.1007/978-3-319-49004-5_31](https://doi.org/10.1007/978-3-319-49004-5_31).

- [44] J. Ramos et al., Using tf-idf to determine word relevance in document queries, in: *Proceedings of the First Instructional Conference on Machine Learning*, Vol. 242, Piscataway, NJ, 2003, pp. 133–142. <https://dl.acm.org/citation.cfm?id=2227115.2227123>.
- [45] M.Q. Raza and A. Khosravi, A review on artificial intelligence based load demand forecasting techniques for smart grid and buildings, *Renewable and Sustainable Energy Reviews* **50** (2015), 1352–1372. doi:[10.1016/j.rser.2015.04.065](https://doi.org/10.1016/j.rser.2015.04.065).
- [46] D. Richards and T. Richardson, Computer-based psychological treatments for depression: A systematic review and meta-analysis, *Clinical psychology review* **32**(4) (2012), 329–342. doi:[10.1016/j.cpr.2012.02.004](https://doi.org/10.1016/j.cpr.2012.02.004).
- [47] G. Rizzo and R. Troncy, Nerd: A framework for unifying named entity recognition and disambiguation extraction tools, in: *Proceedings of the Demonstrations at the 13th Conference of the European Chapter of the Association for Computational Linguistics*, Association for Computational Linguistics, 2012, pp. 73–76. <https://dl.acm.org/citation.cfm?id=2380936>.
- [48] R. Ros, E. Bjarnason and P. Runeson, A machine learning approach for semi-automated search and selection in literature studies, in: *Proceedings of the 21st International Conference on Evaluation and Assessment in Software Engineering, EASE'17*, ACM, 2017, pp. 118–127. ISBN 978-1-4503-4804-1. <https://dl.acm.org/citation.cfm?id=3084243>.
- [49] M. Sabou, D. Winkler, P. Penzerstadler and S. Biffl, Verifying conceptual domain models with human computation: A case study in software engineering, in: *Sixth AAAI Conference on Human Computation and Crowdsourcing*, 2018. <https://www.aaai.org/ocs/index.php/HCOMP/HCOMP18/paper/view/17917/16921>.
- [50] H. Saif, Y. He and H. Alani, Semantic sentiment analysis of Twitter, *The Semantic Web–ISWC* **2012** (2012), 508–524. doi:[10.1007/978-3-642-35176-1_32](https://doi.org/10.1007/978-3-642-35176-1_32).
- [51] A.A. Salatino, F. Osborne, A. Birukou and E. Motta, Improving editorial workflow and metadata quality at springer nature, *International Semantic Web Conference* 2019, p. 2019. <http://oro.open.ac.uk/62265/>.
- [52] A.A. Salatino, F. Osborne and E. Motta, How are topics born? Understanding the research dynamics preceding the emergence of new areas, *PeerJ Computer Science* **3** (2017), e119. doi:[10.7717/peerj-cs.119](https://doi.org/10.7717/peerj-cs.119).
- [53] A.A. Salatino, F. Osborne, T. Thanapalasingam and E. Motta, The CSO classifier: Ontology-driven detection of research topics in scholarly articles, in: *TPDL 2019: 23rd International Conference on Theory and Practice of Digital Libraries*, 2019. doi:[10.1007/978-3-030-30760-8_26](https://doi.org/10.1007/978-3-030-30760-8_26).
- [54] A.A. Salatino, T. Thanapalasingam, A. Mannocci, F. Osborne and E. Motta, The computer science ontology: A large-scale taxonomy of research areas, in: *International Semantic Web Conference*, Springer, 2018, pp. 187–205. doi:[10.1007/978-3-030-00668-6_12](https://doi.org/10.1007/978-3-030-00668-6_12).
- [55] A.A. Salatino, T. Thanapalasingam, A. Mannocci, F. Osborne and E. Motta, Classifying research papers with the computer science ontology, in: *International Semantic Web Conference (P&D/Industry/BlueSky). CEUR Workshop Proceedings*, Vol. 2180, 2018. doi:[10.21954/ou.rd.7204814.v1](https://doi.org/10.21954/ou.rd.7204814.v1).
- [56] M. Sanderson and B. Croft, Deriving concept hierarchies from text, in: *Proceedings of the 22nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, ACM, 1999, pp. 206–213. doi:[10.1145/312624.312679](https://doi.org/10.1145/312624.312679).
- [57] J.M. Schultz and M. Liberman, Topic detection and tracking using idf-weighted cosine coefficient, in: *Proceedings of the DARPA Broadcast News Workshop*, Morgan Kaufmann, San Francisco, 1999, pp. 189–192. <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.591.7692&rep=rep1&type=pdf>.
- [58] A. Sinha, Z. Shen, Y. Song, H. Ma, D. Eide, B.P. Hsu and K. Wang, An overview of Microsoft academic service (MAS) and applications, in: *Proceedings of the 24th International Conference on World Wide Web*, ACM, 2015, pp. 243–246. doi:[10.1145/2740908.2742839](https://doi.org/10.1145/2740908.2742839).
- [59] Y. Sun, Y. Yang, H. Zhang, W. Zhang and Q. Wang, Towards evidence-based ontology for supporting systematic literature review, in: *International Conference on Evaluation and Assessment in Software Engineering (EASE)*, IET, 2012. doi:[10.1049/ic.2012.0022](https://doi.org/10.1049/ic.2012.0022).
- [60] T. Vale, I. Crnkovic, E. Santana de Almeida, P.A. da Mota Silveira Neto, Y. Cerqueira Cavalcanti and S.R. de Lemos Meira, Twenty-eight years of component-based software engineering, *Journal of Systems and Software* **111**(1) (2016), 128–148, ISSN 0164-1212. doi:[10.1016/j.jss.2015.09.019](https://doi.org/10.1016/j.jss.2015.09.019).
- [61] R.J. Wieringa, *Design Science Methodology for Information Systems and Software Engineering*, Springer, Berlin Heidelberg, 2014. doi:[10.1007/978-3-662-43839-8](https://doi.org/10.1007/978-3-662-43839-8).
- [62] M.D. Wilkinson, M. Dumontier, J. Aalbersberg, G. Appleton, M. Axton, A. Baak, N. Blomberg, J.-W. Boiten, L.B. da Silva Santos, P.E. Bourne et al., The fair guiding principles for scientific data management and stewardship, *Scientific data* **3** (2016), 160018. doi:[10.1038/sdata.2016.18](https://doi.org/10.1038/sdata.2016.18).
- [63] G. Wohlgenannt, A. Weichselbraun, A. Scharl and M. Sabou, Dynamic integration of multiple evidence sources for ontology learning, *Journal of Information and Data Management* **3**(3) (2012), 243–254. <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.301.6283&rep=rep1&type=pdf>.
- [64] C. Wohlin and R. Prikladnicki, Systematic literature reviews in software engineering, *Information and Software Technology* **55**(6) (2013), 919–920. doi:[10.1016/j.infsof.2013.02.002](https://doi.org/10.1016/j.infsof.2013.02.002).

- [65] C. Wohlin, P. Runeson, P.A. da Mota Silveira Neto, E. Engström, I. do Carmo Machado and E. Santana de Almeida, On the reliability of mapping studies in software engineering, *The Journal of systems and software* **86**(10) (2013), 2594–2610. doi:[10.1016/j.jss.2013.04.076](https://doi.org/10.1016/j.jss.2013.04.076).
- [66] C. Wohlin, P. Runeson, M. Höst, M.C. Ohlsson, B. Regnell and A. Wesslén, *Experimentation in Software Engineering*, Springer, 2012. doi:[10.1007/978-3-642-29044-2](https://doi.org/10.1007/978-3-642-29044-2).
- [67] N.J.E. Wolfram, P. Lago and F. Osborne, Sustainability in software engineering, in: *IFIP Conference on Sustainable Internet and ICT for Sustainability (SustainIT)*, 2017. doi:[10.23919/sustainit.2017.8379798](https://doi.org/10.23919/sustainit.2017.8379798).
- [68] H. Zhang and M.A. Babar, Systematic reviews in software engineering: An empirical investigation, *Information and Software Technology* **55**(7) (2013), 1341–1354, ISSN 0164-1212. doi:[10.1016/j.infsof.2012.09.008](https://doi.org/10.1016/j.infsof.2012.09.008).
- [69] H. Zhang, M.A. Babar and P. Tell, Identifying relevant studies in software engineering, *Information and Software Technology* **53**(6) (2011), 625–637. doi:[10.1016/j.infsof.2010.12.010](https://doi.org/10.1016/j.infsof.2010.12.010).